

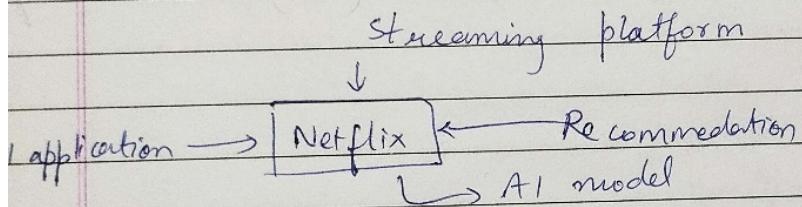
Introduction to Machine Learning

Agenda

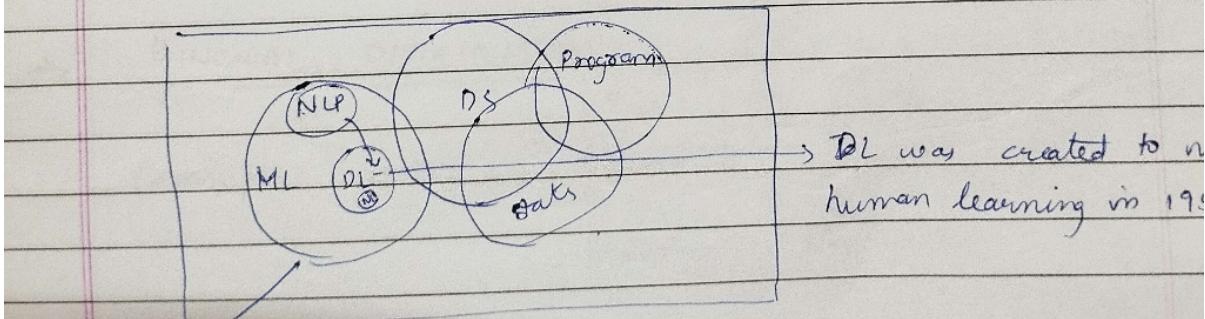
- ① ML Introduction
 - ② AI vs ML vs DL vs DS
 - ③ Simple Linear Regression → Mathematical Intuition behind it.

② AI vs ML vs DL vs DS

AI → creating an application where it performs all its task without human intervention.

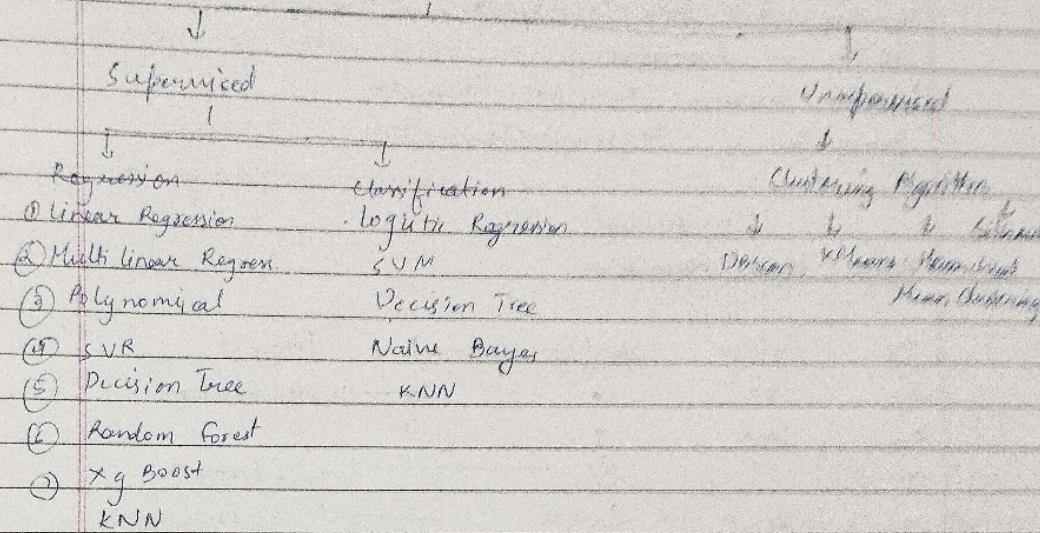


Eg:- Chatbot, self driving car, Auto suggestion, youtube recommendation



^{Provides}
ML → stats tools to analyze, visualize, perform predictions & other tasks with the help of data.

ML & DL



- In Supervised ML → O/P will be given
 → In Unsupervised ML → O/P is not known.

↓
 Here we will make Clusters.

Supervised			Predict	Unsupervised ML
Degree	Exp	Salary	↑ I/O Feature	
B.Tech	2	50k		
PhD	2	70k		
-	-	-	Regression	
-	-	-	! Problem	

No. of play hrs	No. of Study hrs	Independent Feature	Dependent Feature
9	1	(Part/Girl)	O-Offail
7	2	O/Offail	(Part)
3	5	(Part)	

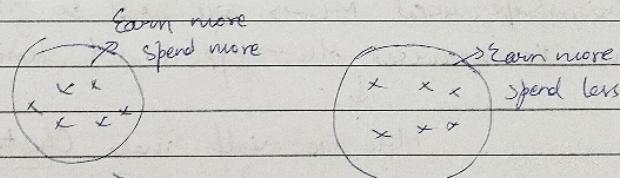
Classification Problem

- Eg:- Flight Price Prediction \rightarrow Regression
 Nigerian Forest Fire \rightarrow Classification
 A&I \rightarrow Regression
 Whether tomorrow Rain/Not \rightarrow Classification
 Buying Day of Person \rightarrow Classification

Unsupervised ML

~~Simple~~ ~~unsupervised~~

Age	Salary	Spending-score (1-10)	Product
24	70k	1	\rightarrow No promotion
26	100k	9	
-	-	2	\rightarrow 20% promotion
21	20k	9	\rightarrow 10% Promotion
25	120k	2	



This particular scenario is "CUSTOMER SEGMENTATION"

SIMPLE LINEAR REGRESSION

\Rightarrow 1 Independent feature 1 Dependent feature -

Eg:- Dataset

HT. WT.

Aim :- If \rightarrow HT, Predict \rightarrow WT

\Rightarrow Eg: (1) Dataset

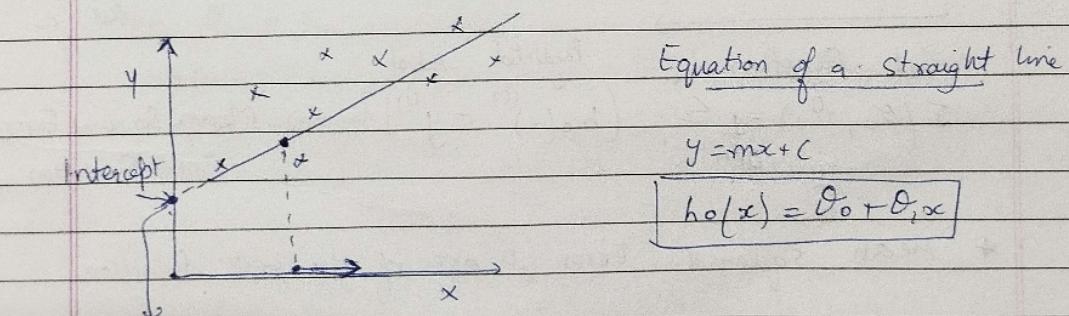
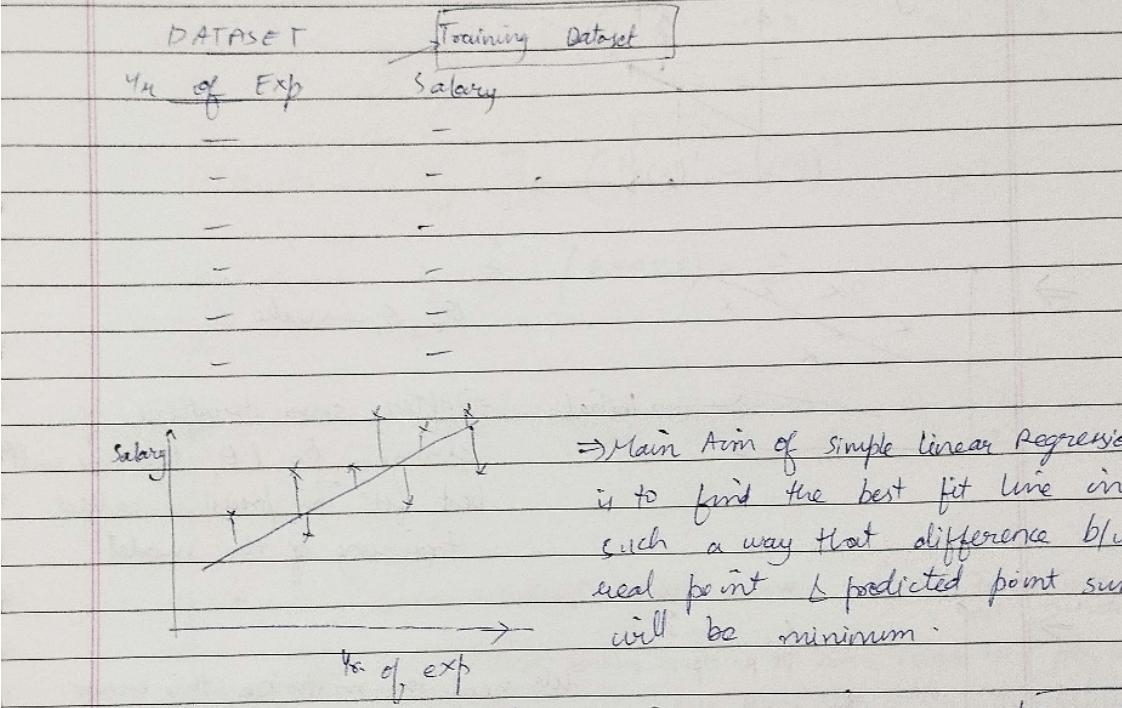
No. of Rooms Price

If No. of Rooms \uparrow Price \uparrow

⇒ Eg :- Yrs of Exp. Salary

Model → will be trained on Yrs. of Exp & salary

Predict → salary based on I/P year

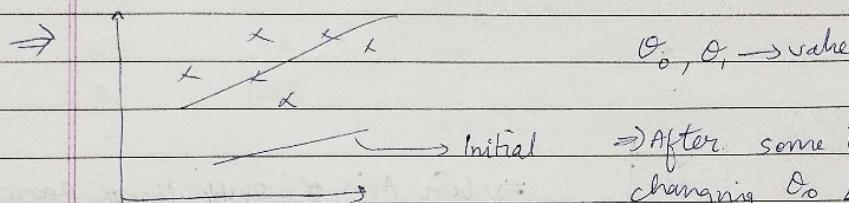
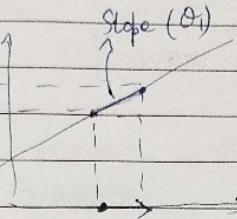


When x value is 0 where is the best fit line meeting y-axis
thus intercept.

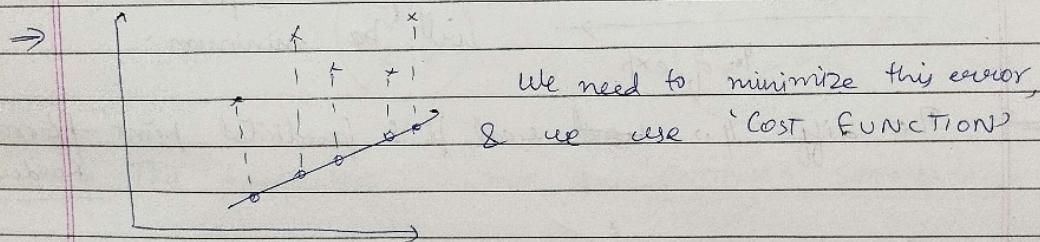
Exp.	Salary
0	3.25 lakh

θ_1 = slope.

Slope \rightarrow With the unit movement in x axis what is movement in y axis



\Rightarrow After some iterations of changing θ_0 & θ_1 , the line will best fit. This process is called training of the model.



$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (\text{Predicted } h_{\theta}(x^{(i)}) - \text{Actual } y^{(i)})^2 \Rightarrow \text{Mean Squared Error}$$

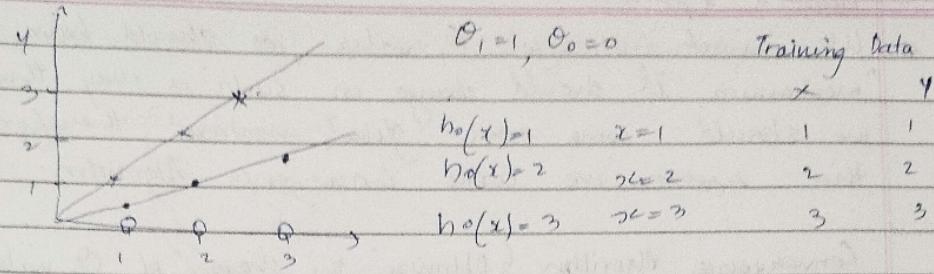
Cost Function

* Mean squared error is one of the cost function

Final Aim \rightarrow in linear Regression

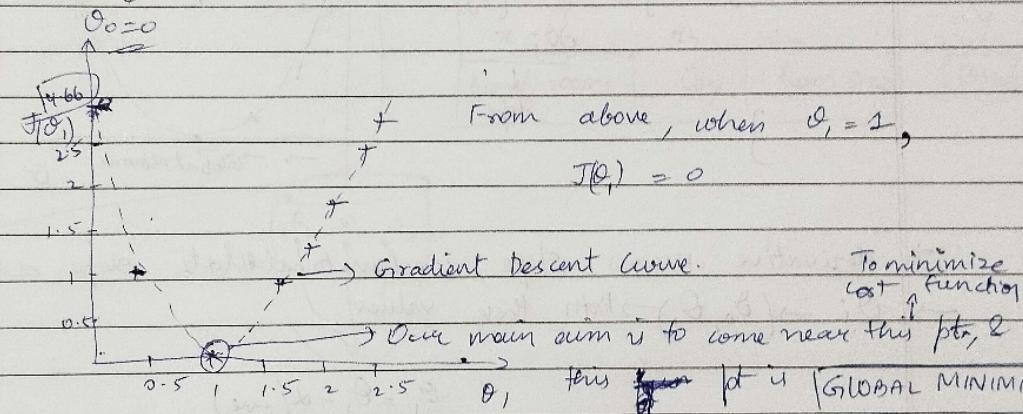
$$\text{Minimize, } J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \underbrace{(h_{\theta}(x^{(i)}) - y^{(i)})^2}_{\text{Mean squared error}}$$

by changing θ_0 & θ_1 values



$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$$J(\theta_1) = \frac{1}{3} \sum_{i=1}^3 (0+0+0) = 0$$



→ When $\theta_1 = 0.5$

$$h_0(x) = 0.5 \quad x = 1$$

$$h_0(x) = 1 \quad x = 2$$

$$h_0(x) = 1.5 \quad x = 3$$

→ Plotting the pts. above.

$$J(\theta_1) = \frac{1}{3} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

$$= \frac{1}{3} [0.25 + 1 + 2.25] = \frac{3.5}{3} = 1.1667$$

→ When $\theta_1 = 0$, → Plotting pts. above

$$h_0(x) = 0$$

$$J(\theta_1) = \frac{1}{3} [(0-1)^2 + (0-2)^2 + (0-3)^2] = \frac{14}{3} = 4.66$$

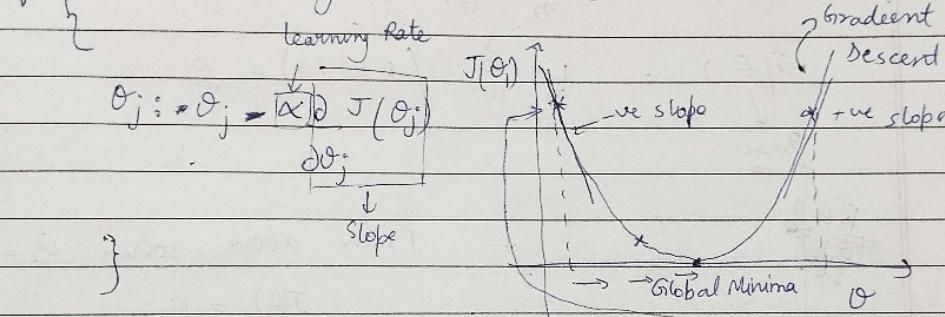
There is a problem

- * We cannot initialize θ_0 value, we should have a mechanism, it should change in such a way that we should come over 'global minima' therefore for this reason we write convergence Algorithm.

Convergence Algorithm {Optimize the changes of θ values}

Repeat until convergence

Also for θ_0 .



- derivative means slope, we have to calculate slope at
- $\theta_j \rightarrow (\theta_0, \theta_j)$ Both these values

$$\theta_j = \theta_j - \alpha (-ve)$$

$$= \theta_j + \alpha \rightarrow \text{Increasing } \theta_j$$

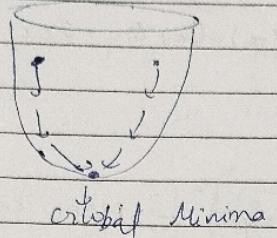
- When we plot this line downwards we have to increase θ_j to get it reach Global minima.
- Slope is used to update is (θ_0, θ_j) value.

- In case of that 2nd the slope to reach global minima we have to subtract α

$$\boxed{\theta_j = \theta_j - \alpha (+ve slope)}$$

- $\alpha \rightarrow$ Learning Rate \rightarrow How fast you want to converge.
Usually selected $\boxed{\alpha = 0.001}$

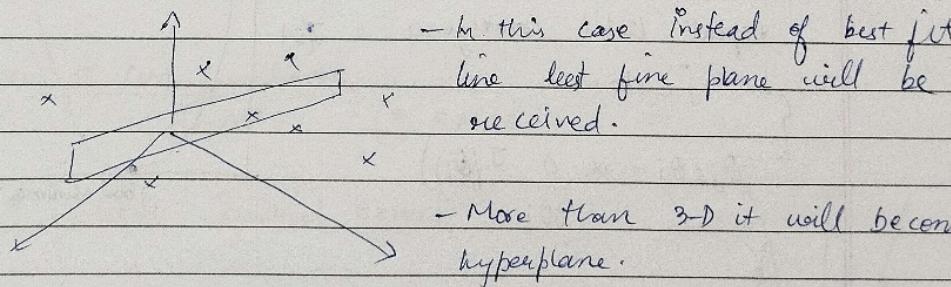
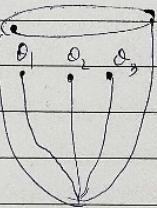
→ If we include θ_0 also, we will get a 3-D shape -



for Multilinear Regression

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

x_1	x_2	x_3	y
No. of rooms	City	Room size	Prize



Linear Regression Algorithm

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_n x_n$$

Convergence Algorithm

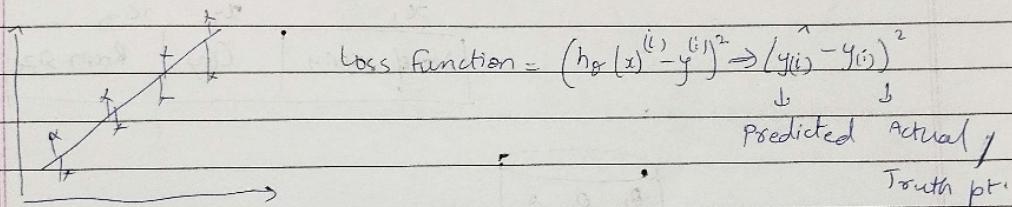
Cost function,

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad \text{MEAN SQUARED ERROR}$$

\Rightarrow [Loss function vs Cost function]

One Data point

All observation

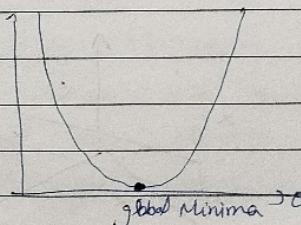


CONVERGENCE THEOREM

repeat until convergence

$j = 1 \dots m$

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_j)}{\partial \theta_j}$$



$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_0} \left[\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right]$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$J = \frac{\partial}{\partial \theta_0} \left[\frac{1}{2m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})^2 \right]$$

$$= \frac{1}{m} \sum_{i=1}^m \left[(\theta_0 + \theta_1 x_i) - y^{(i)} \right]^2$$

$$\text{J}_{\theta} = \frac{\partial}{\partial \theta_1} \left[\frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)}) - y^{(i)} \right]^2$$

$$= \frac{1}{m} \sum_{i=1}^m \left[(\theta_0 + \theta_1 x^{(i)}) - y^{(i)} \right] [x]$$

Repeat Until Convergence

$$\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_0(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_1(x^{(i)}) - y^{(i)}) \end{array} \right.$$

α -Speed of convergence
Global Minim

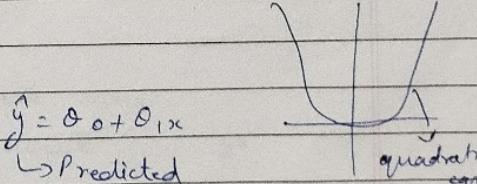
Types of Cost functions.

- ① MSE
- ② MAE
- ③ RMSE

① MSE { Mean Squared Error }

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

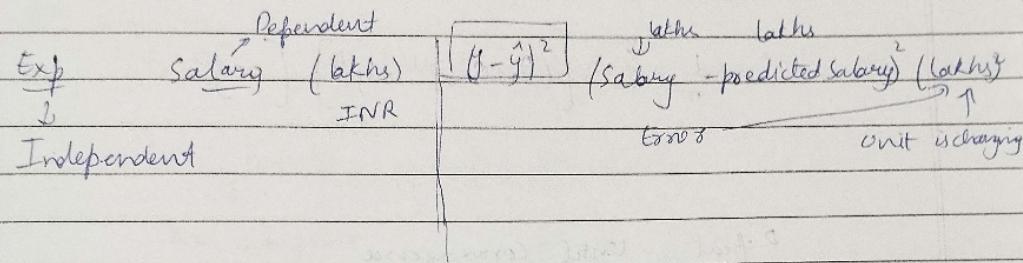
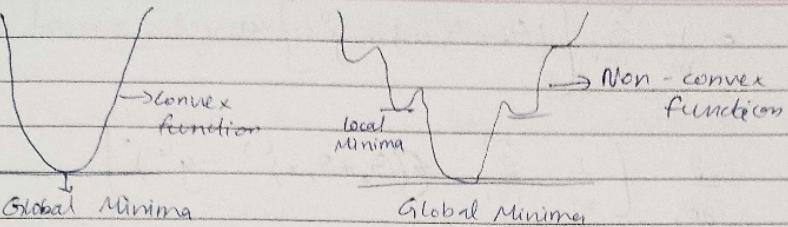
$\hat{y} = \theta_0 + \theta_1 x$
↳ Predicted



Advantages :-

Disadvantage

- i) This equation is differentiable
- ii) This equation also has one global minima.
- i) Not Robust to outliers.
- ii) Penalizing the error.



② MAE (Mean Absolute Error)

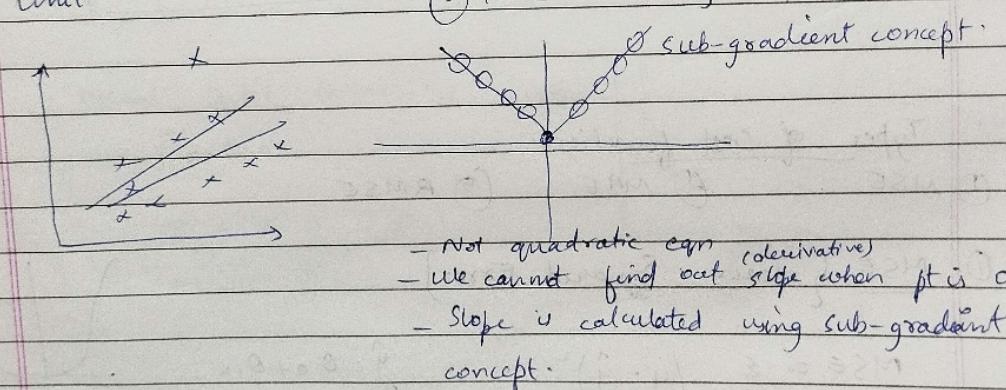
$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y - \hat{y}|$$

Advantage

- ① Robust to outliers
- ② It will also be same unit

Disadvantage

- ① Convergence usually takes more time
- ② Optimization is a complex task.
- ③ Time consuming.



Advantages & Disadvantages

(1) Huber Loss



MSE & MAE

Unit, Outliers, Differentiable

(2) R MSE



\sqrt{MSE}

RANKA
Date / /
Page / /

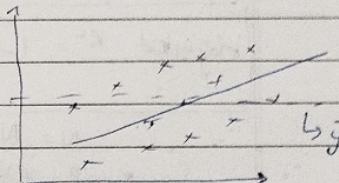
Performance Matrix

(1) R Squared

(2) Adjusted R squared

(3) R Squared

$$R^{\text{Squared}} = 1 - \frac{SS_{\text{Res}}}{SS_{\text{Total}}} = + -$$



S_{Res} = Sum of square Residuals

S_{Total} = Sum of Square Average

$$\bar{y} = \text{Avg. of } y$$

$$= 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \rightarrow \text{if model is fitted well it will be 1}$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 \rightarrow \text{high}$$

$$R^{\text{Squared}} = 1 - \left\{ \begin{array}{l} \text{Small number} \\ \text{Bigger number} \end{array} \right\} \rightarrow \text{small}$$

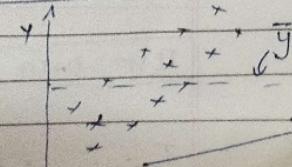
R-squared / of ≤ 1

$\rightarrow 0.85 \rightarrow 85\% \text{ Accurate Model}$

$> 0.75 \rightarrow 75\% \text{ Accurate}$

R-square checks performance of the model created

If R-square is -ve \rightarrow very bad model.



* Adjusted R squared

Size of House | City location | No. of bedrooms | Gender | Price

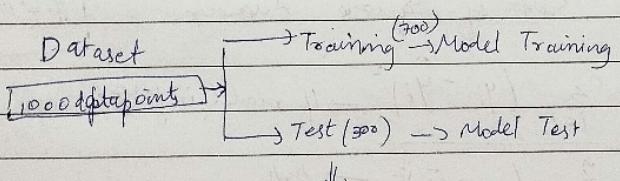
R^2 65% Adj R^2 = 63%, $P=1$
 R^2 75% Adj R^2 = 73%, $P=2$
 R^2 88%
 R^2 90% Should NOT have happened
 \rightarrow Adj R^2 85%

$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2)(N-1)}{N-P-1}$$

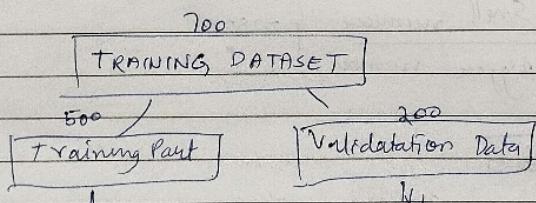
N = No. of data points

P = No. of independent features

Overfitting and Underfitting (Bias and Variance)



New Dataset



Training the model

Hyperparameters tuning the model.

TRAIN DATA	Very Good Accuracy (90%) [low Bias]
------------	-------------------------------------

TEST DATA	Very Good Accuracy (85%) [low variance] ↳ Generalized Model
-----------	--

Training data info is passed by Bias, test data using variance.

Train	Bad Accuracy (High Bias)
Test	Bad Accuracy (High Variance)

Eg (1) Train → Very Good Accuracy (90%) → [low Bias] → Solve this problem by
 Test → Bad Accuracy (50%) → [high Variance] → hyperparameter tuning
 ↓
 Overfitting

(2) Train → [High Bias]
 Test → Model is low/High → [low or high variance]

↓
 Underfitting

